

## LayerNorm vs RMSNorm

Original transformer: **LayerNorm** – normalizes the mean and variance across  $d_{model}$

$$y = \frac{x - \mathbf{E}[x]}{\sqrt{\mathbf{Var}[x] + \epsilon}} * \gamma + \beta$$

Many modern LMs: **RMSNorm** – does not subtract mean or add a bias term

$$y = \frac{x}{\sqrt{\|x\|_2^2 + \epsilon}} * \gamma$$

**Notable models:**

GPT3/2/1, OPT, GPT-J, BLOOM

**Notable models:**

LLaMA-family, PaLM, Chinchilla, T5